

# Self-healing of operational workflow incidents on distributed computing infrastructures

Rafael Ferreira da Silva\*, Tristan Glatard\*, Frédéric Desprez†

\*University of Lyon, CNRS, INSERM, CREATIS, Villeurbanne, France

†INRIA, University of Lyon, LIP, ENS Lyon, Lyon, France

Email: {rafael.silva,glatard}@creatis.insa-lyon.fr, Frederic.Desprez@inria.fr o

**Abstract**—Distributed computing infrastructures are commonly used through scientific gateways, but operating these gateways requires important human intervention to handle operational incidents. This paper presents a self-healing process that quantifies incident degrees of workflow activities from metrics measuring long-tail effect, application efficiency, data transfer issues, and site-specific problems. These metrics are simple enough to be computed online and they make little assumptions on the application or resource characteristics. Incidents are classified in levels and associated to sets of healing actions that are selected based on association rules modeling correlations between incident levels. The healing process is parametrized on real application traces acquired in production on the European Grid Infrastructure. Implementation and experimental results obtained in the Virtual Imaging Platform show that the proposed method speeds up execution up to a factor of 4 and properly detects unrecoverable errors.

## I. INTRODUCTION

Distributed computing infrastructures (DCI) are becoming daily instruments of scientific research, in particular through scientific gateways [1] developed to allow scientists to transparently run their analyses on large sets of computing resources. While these platforms provide important amounts of resources in an almost seamless way, their large scale and the number of middleware systems involved lead to many errors and faults. Easy-to-use interfaces provided by these gateways exacerbate the need for properly solving operational incidents encountered on DCIs since end users expect high reliability and performance with no extra monitoring or parametrization from their side. In practice, such services are often backed by substantial support staff who monitors running experiments by performing simple yet crucial actions such as rescheduling tasks, restarting services, killing misbehaving experiments or replicating data files to reliable storage facilities. Fair QoS can then be delivered, yet with important human intervention.

For instance, the long-tail effect [2] is a common frustration for users who have to wait for a long time to retrieve the last few pieces of their computations. Operators may be able to address it by rescheduling tasks that are considered late (e.g. due to execution on a slow machine, low network throughput or just loss of contact) but detection is very time consuming and still rough.

Automating such operations is challenging for two reasons. First, the problem is online by nature because no reliable user activity prediction can be assumed, and new workloads may arrive at any time. Therefore the considered metrics, decisions

and actions have to remain simple and to yield results while the application is still executing. Second, it is non-clairvoyant due to the lack of information about applications and resources in production conditions. Computing resources are usually dynamically provisioned from heterogeneous clusters, clouds or desktop grids without any reliable estimate of their availability and characteristics. Models of application execution times are hardly available either, in particular on heterogeneous computing resources.

A scientific gateway is considered here as a platform where users can process their own data with predefined applications workflows. Workflows are compositions of *activities* defined independently from the processed data and that only consist of a program description. At runtime, activities receive data and spawn invocations from their input parameter sets. Invocations are assumed independent from each other (bag of tasks) and executed on the DCI as single-core tasks which can be resubmitted in case of failures. This model fits several existing gateways such as e-bioinfra [3], P-Grade [4], the Virtual Imaging Platform [5] or Décryphon [6]. We also consider that files involved in workflow executions are accessed through a single file catalog but storage is distributed. Files may be replicated to improve availability and reduce load on servers.

The gateway may take decisions on file replication, resource provisioning, and task scheduling on behalf of the user. Performance optimization is a target but the main point is to ensure that correctly-defined executions complete, that performance is acceptable, and that misbehaving runs (e.g. failures coming from user errors or unrecoverable infrastructure downtimes) are quickly detected and stopped before they consume too many resources.

Our ultimate goal is to reach a general model of such a scientific gateway that could autonomously detect and handle operational incidents. In this work, we propose a healing process for workflow activities only. Activities are modeled as Fuzzy Finite State Machines (FuSM) [7] where state degrees of membership are determined by an external healing process. Degrees of membership are computed from metrics assuming that incidents have outlier performance, e.g. a site or a particular invocation behaves differently than the others. Based on incident degrees, the healing process determines incident levels from thresholds determined from platform history. A specific set of actions is then selected from association rules among incident levels.

Section II presents related work. Our approach is described

in section III (general healing process), section IV (metrics used to quantify incident degrees) and section V (incident levels and associated action sets). Experimental results are presented in section VI in production conditions.

## II. RELATED WORK

Managing systems with limited intervention of system administrators is the goal of autonomic computing [8]. It has been used to address various problems related to self-healing, self-configuration, self-optimization, and self-protection of distributed systems. For instance, provisioning of virtual machines is studied by Nguyen et al. [9] and an approach to tackle service overload, queue starvation, “black hole” effect and job failures is sketched by Collet et al. [10].

An autonomic manager consists of monitoring, analysis, planning, execution and knowledge (so-called MAPE-K loop). Generic software frameworks have been built to wrap legacy applications in such loops with limited intrusiveness. For instance, Broto et al. [11] demonstrates the wrapping of DIET grid services for autonomic deployment and configuration. We consider here that the target gateway can be instrumented to report appropriate events and perform actions.

Monitoring is broadly studied in distributed systems, both at coarse (traces, archives) and fine time scales (active monitoring, probing). Many workload archives are available. In particular, the grid observatory [12] has been collecting traces for a few years on several grids. However, as noted by Iosup & Epema [13], most existing traces remain at the task level and lack information about workflows and activities. Application patterns can be retrieved from logs (e.g. bag of tasks) but precise information about workflow activities is bound to be missing. Studies on task errors and their distributions are also available [14], [15], but they do not consider operational issues encountered by the gateways submitting these tasks. Besides, active monitoring using tools such as Nagios [16] cannot be the only monitoring source when substantial workloads are involved. Therefore we rely on traces of the target gateway, as detailed in section V. One issue in this case is to determine the timespan where system behavior can be considered steady-state. Although this issue was recently investigated [17], it remains difficult to identify non-stationarities in an online process and we adopt here a stationary model.

Analysis consists in computing metrics (a.k.a. utility functions) from monitoring data to characterize the state of the system. System state usually distinguishes two regimes: properly functioning and malfunctioning. Zhang et al. [18] assume that incidents lead to non stationarity of the workload statistics and use the Page-Hinkely test to detect them. Stehle et al. [19] present a method where the convex hull is used instead of hyper-rectangles to classify system states. As described in section V, we use multiple threshold values for a given metric to use more than two levels to characterize incidents.

Planning and actions considered in this work deal with task scheduling and file replication. In these domains, most approaches are clairvoyant, meaning that resource, task, error rate and workload characteristics are precisely known [20], [21]. Heuristics are designed by Casanova et al. [22] for the

case where only data transfer costs are known, on an offline problem though. Camarasu-Pop et al. [23] propose a dynamic load-balancing strategy proposed to remove the long-tail effect on production heterogeneous systems, but it is limited to Monte-Carlo simulations.

The general task scheduling problem is out of our scope. We assume that a scheduler is already in place, and we only aim at performing actions when it does not deliver proper performance. In particular, we focus on site blacklisting and on task replication to avoid long-tail effect. Task replication, a.k.a. redundant requests is commonly used to address non-clairvoyant problems [2], but it should be used sparingly to avoid overloading the middleware and degrading fairness among users [24]. In this work, task replication is considered when activities are detected blocked or of low efficiency according to the metric presented in section IV.

Similarly, file replication strategies often assume clairvoyance on the size of produced data, file access pattern and infrastructure parameters [25], [26]. In practice, production systems mostly remain limited to manual replication strategies though [27].

## III. GENERAL HEALING PROCESS

An activity is modeled as an FuSM with 13 states shown on Fig. 1. The activity is initialized in *Submitting Invocations* where all the tasks are generated and submitted. Tasks consist of 4 successive phases: initialization, inputs download, application execution and output upload. They are all assumed independent, but with similar execution times (bag of tasks). *Running* is a state where no particular issue is detected; no action is taken and the activity is assumed to behave normally. *Completed* (resp. *Failed*) is a terminal state used when all the invocations are successfully completed (resp. at least one invocation failed). These 4 states are crisp (not fuzzy) and exclusive<sup>1</sup>. The 9 other states are fuzzy states corresponding to detected incidents.

The healing process sets the degree of FuSM states from incident detection metrics and invocation statuses. Then it determines the actions to be performed to address the incidents. If no action is required then the process waits until an event occurs (task status change) or a timeout is reached.

Let  $I = \{x_i, i = 1, \dots, n\}$  be the set of possible incidents (9 in this work) and  $\eta = (\eta_1, \dots, \eta_n) \in [0, 1]^n$  their degrees in the FuSM. Incident  $x_i$  can occur at different levels  $\{x_{i,j}, j = 1, \dots, m_i\}$  delimited by threshold values  $\tau_i = \{\tau_{i,j}, j = 1, \dots, m_i\}$ . The level of incident  $i$  is determined by  $j$  such that  $\tau_{i,j} \leq \eta_i < \tau_{i,j+1}$ . A set of actions  $a_i(j)$  is available to address  $x_{i,j}$ :

$$\begin{aligned} a_i : [1, m_i] &\rightarrow \wp(A) \\ j &\mapsto a_i(j) \end{aligned} \quad (1)$$

where  $A$  is the set of possible actions taken by the healing process and  $\wp(A)$  is the power set of  $A$ .

In addition to the incidents themselves, incident causes are taken into account. Association rules [28] are used to identify

<sup>1</sup>Their degree can only be 0 or 1 and if 1 then all the other states have a degree of 0.

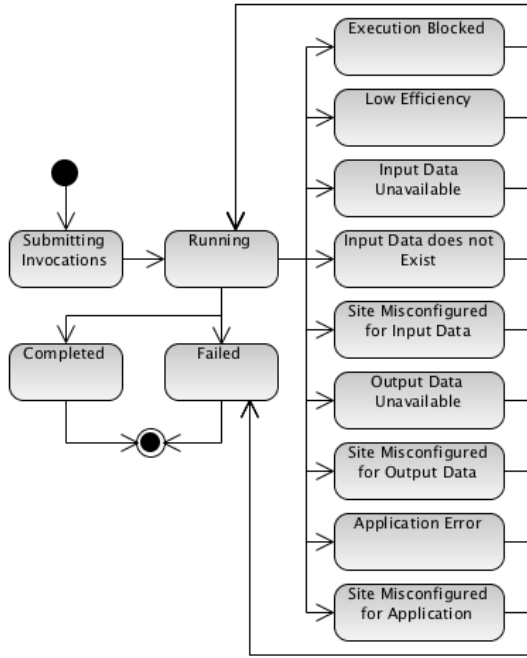


Fig. 1. Fuzzy Finite State Machine (FuSM) representing an activity.

relations between levels of different incidents. Association rules to  $x_{i,j}$  are defined as  $R_{i,j} = \{r_{i,j}^{u,v} = (x_{u,v}, x_{i,j}, \rho_{i,j}^{u,v})\}$ . Rule  $r_{i,j}^{u,v}$  means that when  $x_{u,v}$  happens then  $x_{i,j}$  also happens with confidence  $\rho_{i,j}^{u,v} \in [0, 1]$ . The confidence of a rule is an estimate of probability  $P(x_{i,j}|x_{u,v})$ . Note that  $r_{i,j}^{i,j} \in R_{i,j}$  and  $\rho_{i,j}^{i,j} = 1$ . We also define  $R = \bigcup_{i \in [1, n], j \in [1, m_i]} R_{i,j}$ .

Fig. 2 presents the algorithm used at each iteration of the healing process. Incident degrees are determined based on metrics presented in section IV and incident levels  $j$  are obtained from historical data as explained in section V. A roulette wheel selection [29] based on  $\eta$  is then performed to select  $x_{i,j}$  the incident level of interest at this iteration. Roulette wheel selection assigns a proportion of the wheel to each incident level according to their probability and a random selection is performed based on a spin of the roulette wheel. The probability of an incident  $x_i$  to be selected is  $p(x_i) = \eta_i / \sum_{j=1}^n \eta_j$ . A potential cause  $x_{u,v}$  for incident  $x_{i,j}$  is then selected from a roulette wheel selection on the association rules  $r_{i,j}^{u,v}$ , where  $x_u$  is at level  $v$ . Rule  $r_{i,j}^{u,v}$  is weighted  $\eta_u \times \rho_{i,j}^{u,v}$  in the roulette selection. Only first-order causes are considered here but the approach could be extended to include more recursion levels. Note that  $r_{i,j}^{i,j}$  participates in this selection so that a first-order cause is not systematically chosen. Finally, actions in  $a_u(v)$  are performed.

Table I illustrates this mechanism on an example case where only 3 incidents are considered.

#### IV. INCIDENT DEGREE

This section describes the metrics used to determine the degree of the 9 considered incidents (step 02 on Fig. 2).

a) *Activity Blocked*: this incident happens when an invocation is considered late compared to the others. It is responsible for many operational issues, leading to substantial speed-up reductions. For instance, it occurs when one invocation of the

**Input:** invocation statuses and history of  $\eta$

**Output:** set of actions  $a$

01. wait for event or timeout
02. determine incident degrees  $\eta$  based on metrics
03. determine incident levels  $j$  such that  $\tau_{i,j} \leq \eta_i < \tau_{i,j+1}$
04. select incident  $x_i$  by roulette wheel selection based on  $\eta$
05. select rule  $r_{u,v} = (x_{u,v}, x_{i,j}, \rho_{i,j}^{u,v}) \in R_{i,j}$  by roulette wheel selection based on  $\eta_u \times \rho_{i,j}^{u,v}$ , where  $x_u$  is at level  $v$
06.  $a = a_u(v)$
07. perform actions in  $a$

Fig. 2. One iteration of the healing process.

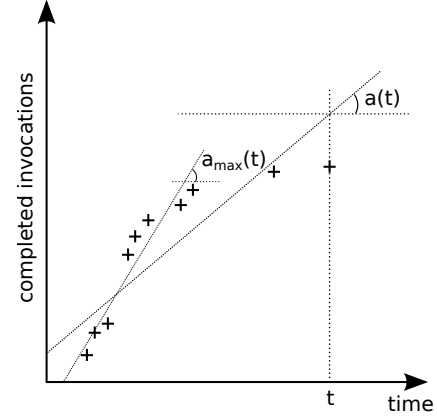


Fig. 3. Detection of blocked activity.

activity requires more CPU cycles or when the invocation faces longer waiting times, lost tasks or executes on resources with poorer performance. This situation is detected online from the number  $n(t)$  of completed invocations at time  $t$  (see Fig. 3). At time  $t$ , we compute the slope  $a(t)$  of the regression line of  $\{(t_i, n(t_i)), t_i \leq t\}$ . In case the iteration is triggered by a timeout instead of an event, then  $(t, n(t) + 1)$  is added to the regression set. This is meant to ensure that long-running invocations can be handled before they complete. We then define the incident degree  $\eta_b$  from the contraction rate of the

$x_i$ : incident name	Degree $\eta_i$	Level $j$
$x_1$ : activity blocked	0.8	2
$x_2$ : low efficiency	0.4	1
$x_3$ : input data unavailable	0.1	1

Step 04:  $x_{1,2}$  is selected with probability  $\frac{0.8}{0.8+0.4+0.1}$ .

Rule	Confidence
$r_{1,2}^{2,1}: x_{2,1} \rightarrow x_{1,2}$	0.8
$r_{1,2}^{3,1}: x_{3,1} \rightarrow x_{1,2}$	0.2
$r_{1,2}^{1,2}: x_{1,2} \rightarrow x_{1,2}$	1

$r_{1,2}^{2,1}$  is chosen with probability  $\frac{0.8 \times 0.4}{0.8 \times 0.4 + 0.2 \times 0.1 + 0.8 \times 1}$ .

Step 06: actions in  $a_2(1)$  are performed.

TABLE I  
EXAMPLE CASE.

linear regression slope:

$$\eta_b = 1 - \frac{a(t)}{a_{\max}(t)}$$

where  $a_{\max}(t)$  is maximal value of  $a(t)$  in  $[0, t]$ .  $t = 0$  is the time when the activity is started, i.e., all the invocations are initialized. Note that the maximum degree  $\eta_b = 1$  is reached when the activity is completely blocked ( $\lim_{t \rightarrow \infty} a(t) = 0$ ). On the other hand,  $\eta_b = 0$  is reached when  $a(t) = a_{\max}(t)$ . Invocations are assumed of identical lengths, which is common for a workflow activity.

*b) Low Efficiency:* this happens when the time spent by all the activity invocations in data transfers dominates CPU time. It may be due to sites with poor network connectivity or intrinsic to the application. The incident degree is defined from the ratio between the cumulative CPU time  $C_i$  consumed at time  $t$  by all completed invocations and the cumulative execution time at time  $t$  of all completed invocations:

$$\eta_e = 1 - \frac{\sum_{i=1}^{n(t)} C_i}{\sum_{i=1}^{n(t)} (C_i + D_i)}$$

where  $D_i$  is the time spent by invocation  $i$  in data transfers.

*c) Input Data Unavailable:* this happens when a file is registered in the file catalog but the storage resource(s) is(are) unavailable or unreachable. The incident degree  $\eta_{iu}$  in this state is determined from the input transfer failure rate due to data unavailability. Transfers of completed, failed, and running invocations are considered.

*d) Input Data does not Exist:* this happens when an incorrect data path was specified, the file was removed by mistake or the file catalog is unavailable or unreachable. Again, the incident degree  $\eta_{ie}$  is directly determined by the input transfer failure rate due to non-existent data. Transfers of completed, failed, and running invocations are considered.

*e) Site Misconfigured for Input Data:* this incident happens when sites have utmost input data transfer failure rate. The incident degree  $\eta_{is}$  at time  $t$  is measured as follows:

$$\eta_{is} = \max(\phi_1, \phi_2, \dots, \phi_k) - \text{median}(\phi_1, \phi_2, \dots, \phi_k)$$

where  $\phi_i$  denotes the input transfer failure ratio (including both input data unavailable and input data does not exist) on site  $i$  at time  $t$  and  $k$  is the number of white-listed sites used by the activity at time  $t$ . The difference between the maximum rate and the median ensures that the incident degree has high values only when some sites are misconfigured. This metric is correlated but not redundant with the two previous ones. If some input data file is not available due to site-independent issues with the storage system, then  $\eta_{iu}$  will grow but  $\eta_{is}$  will remain low because all sites fail identically. On the contrary,  $\eta_{is}$  may grow while  $\eta_{iu}$  and  $\eta_{ie}$  remain low.

*f) Output Data Unavailable:* output data can also be unavailable. Unavailability happens due to three main reasons: the user did not specify the output path correctly, the application did not produce the expected data, or the file catalog or storage resource are unavailable or unreachable. The incident degree  $\eta_{ou}$  is determined by the output transfer failure rate. Transfers of completed, failed and running invocations are considered.

*g) Site Misconfigured for Output Data:* the incident degree  $\eta_{os}$  in this incident is determined as follows:

$$\eta_{os} = \max(\psi_1, \psi_2, \dots, \psi_k) - \text{median}(\psi_1, \psi_2, \dots, \psi_k)$$

where  $\psi_i$  denotes the output transfer failure ratio on site  $i$  at time  $t$  and  $k$  is the number of white-listed sites used by the activity at time  $t$ .

*h) Application Error:* applications can fail due to a variety of reasons among which: the application executable is corrupted, dependencies are missing, or the executable is not compatible with the execution host. The incident degree  $\eta_a$  in this state is measured by the task failure rate due to application errors. Completed, failed, and running tasks are considered.

*i) Site Misconfigured for Application:* The incident degree  $\eta_{as}$  in this state is measured as follows:

$$\eta_{as} = \max(\alpha_1, \alpha_2, \dots, \alpha_k) - \text{median}(\alpha_1, \alpha_2, \dots, \alpha_k)$$

where  $\alpha_i$  denotes the task failure rate due to application errors on site  $i$  and  $k$  is the number of white-listed sites used by the activity at time  $t$ .

## V. INCIDENT LEVELS AND ACTIONS

Incident degrees  $\eta_i$  are quantified in discrete incident levels so that different sets of actions can be used to address different levels of the incident. The threshold number and values are determined from observed distributions of  $\eta_i$ . The number  $m_i$  of incident levels associated to incident  $i$  is set as the number of modes in the distribution of  $\eta_i$ . Thresholds  $\tau_{i,j}$  are determined from mode clustering.

### A. Training Dataset

We collected traces from the Virtual Imaging Platform [5] gateway between April and August 2011. Applications deployed in this platform are described as workflows executed using the MOTEUR workflow engine [30]. Resource provisioning and task scheduling is provided by DIRAC [31] using so-called ‘‘pilot jobs’’. Resources are provisioned online with no advance reservations. Tasks are executed on the biomed virtual organization (VO) of the European Grid Infrastructure (EGI)<sup>2</sup> which has access to some 150 computing sites worldwide and to 120 storage sites providing approximately 4 PB of disk.

This data set contains 1,082 executions of 36 different workflows executed by 26 users. Workflow executions contain 1,838 activity instances, corresponding to 92,309 invocations and 123,025 tasks (including resubmissions).

Figure 4 shows the cumulative amount of running activities along this period. It shows that the workload is quite uniformly distributed although a slight increase is observed in June.

<sup>2</sup><http://www.egi.eu>

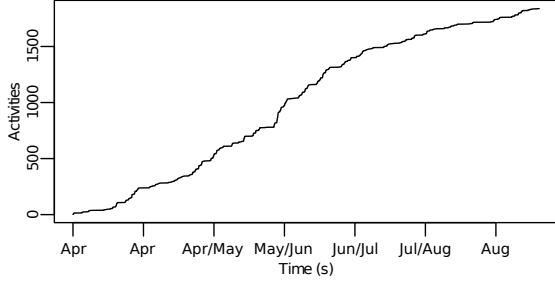


Fig. 4. Cumulative amount of running activities from April to August 2011.

### B. Incident Levels and Actions

Incident degrees were computed after each event found in this data set (total of 641,297 events). Fig. 5 displays histograms of computed incident degrees. For readability purposes, only  $\eta_i \neq 0$  values are represented. Histograms are clearly multi-modal, which confirms that incident degrees are quantified. Level numbers and threshold values  $\tau$  are set from visual mode detection in these histograms and reported on Table II with associated actions.

Incidents at level 1 are considered painless for the execution and they do not trigger any action. Other levels can lead to radical (completely stop the activity or blacklist a site) or intermediate actions (task replication, file replication, or provisioning of extra resources).

### C. Association Rules

Association rules are computed based on the frequency of occurrences of two incident levels. The confidence  $\rho_{i,j}^{u,v}$  of a rule  $x_{u,v} \Rightarrow x_{i,j}$  measures the probability that an incident level  $x_{i,j}$  happens when  $x_{u,v}$  occurs. Table III shows rule samples extracted from the training data-set and ordered by decreasing confidence. The set of rules leading to activity blocked ( $x_{1,2}$ ) and low efficiency ( $x_{2,2}$ ) incidents shows that they are partially dependent of other “cause” incidents, which is considered by the self-healing process.

At the bottom of the table we find rules with null confidence. These are consistent with common-sense interpretation of the incident dependencies (e.g. no site-specific issue when input data is unavailable).

## VI. EXPERIMENTS

The healing process was implemented in the Virtual Imaging Platform (see description in section V-A) and deployed in production. The experiments presented hereafter evaluate the ability of the healing process to (i) improve workflow makespan in case of recoverable incidents and (ii) quickly identify and report critical issues.

### A. Implementation

The FuSM and healing process were implemented in the MOTEUR workflow engine. The timeout value in the healing process was computed dynamically as the median of the task inter-completion delays in the current execution.

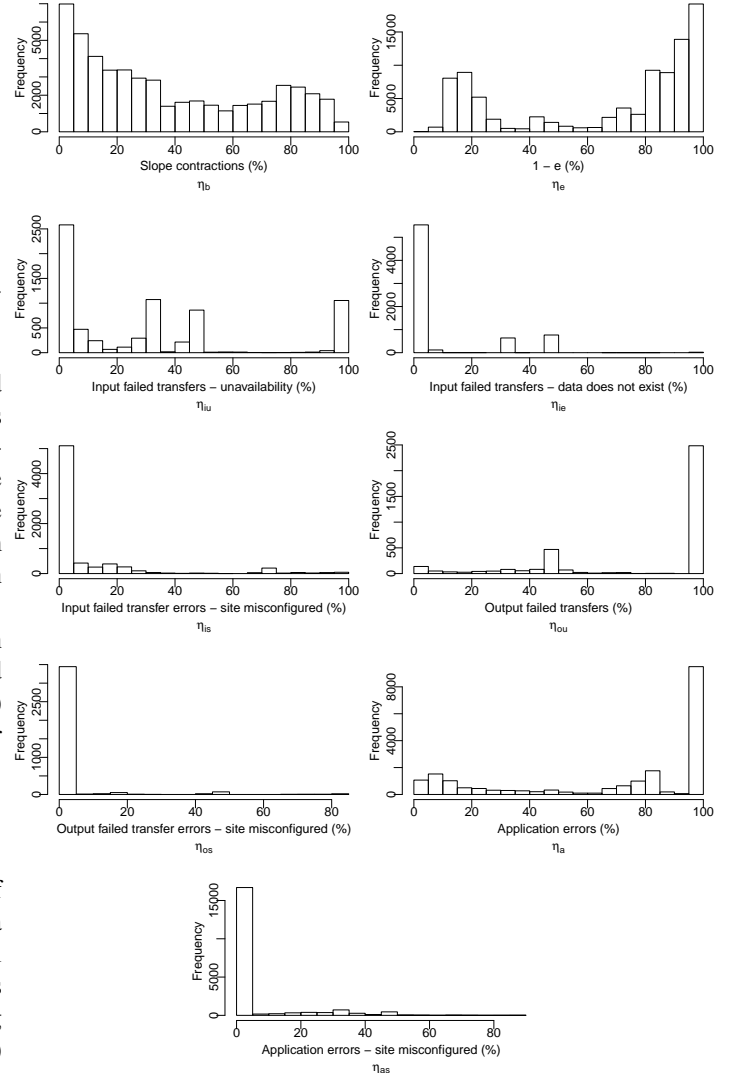


Fig. 5. Histograms of incident degrees sampled in bins of 5%.

Task replication is performed by resubmitting running tasks to DIRAC. To avoid concurrency issues in the writing of output files, a simple mechanism based on file renaming was implemented. To limit infrastructure overload, running tasks are replicated up to 5 times.

Input file unavailability is distinguished from non-existent file using ad-hoc parsing of standard error files. File replication is implemented differently depending on the incident. In case of input data unavailability, a file is replicated to a storage resource randomly selected in the biomed VO. The maximal allowed number of file replicas is set to 5. In case a site is misconfigured, replication to the site local storage resource is first attempted. This aims at circumventing inter-domain connectivity issues. If there is no local storage available or the replication process fails, then a second attempt is performed to a storage resource successfully accessed by other tasks executed on the same site.

Problematic sites are only temporarily blacklisted during a time interval set from exponential back-off. The site is

Incident	Number of incident levels ( $m_i$ )	Level 1		$\tau_{i,2}$	Level 2 actions	Level 3	
		$\tau_{i,1}$	actions			$\tau_{i,3}$	actions
$x_1$ : activity blocked	2	0	$\emptyset$	0.6	replicate running tasks		
$x_2$ : low efficiency	2	0	$\emptyset$	0.6	replicate input files replicate running tasks		
$x_3$ : input data unavailable	3	0	$\emptyset$	0.2	replicate input files	0.8	stop activity
$x_4$ : input data does not exist	2	0	$\emptyset$	0.8	stop activity		
$x_5$ : site misconfigured for input data	3	0	$\emptyset$	0.3	replicate files on sites reachable from problematic site	0.65	blacklist site
$x_6$ : output data unavailable	2	0	$\emptyset$	0.8	stop activity		
$x_7$ : site misconfigured for output data	2	0	$\emptyset$	0.1	blacklist site		
$x_8$ : application error	2	0	$\emptyset$	0.5	stop activity		
$x_9$ : site misconfigured for application	2	0	$\emptyset$	0.1	blacklist site		

TABLE II  
INCIDENT LEVELS AND ACTIONS.

Association rule	$\rho_{i,j}^{u,v}$
$x_{5,2} \Rightarrow x_{2,2}$	0.3809
$x_{7,2} \Rightarrow x_{1,2}$	0.3529
$x_{5,3} \Rightarrow x_{1,2}$	0.3333
$x_{1,2} \Rightarrow x_{2,2}$	0.3059
$x_{3,2} \Rightarrow x_{1,2}$	0.2975
$x_{7,2} \Rightarrow x_{2,2}$	0.2941
$x_{5,2} \Rightarrow x_{1,2}$	0.2608
$x_{9,2} \Rightarrow x_{1,2}$	0.2435
$x_{2,2} \Rightarrow x_{1,2}$	0.2383
...	...
$x_{3,2} \Rightarrow x_{2,2}$	0.1276
$x_{7,2} \Rightarrow x_{3,3}$	0.1250
$x_{3,3} \Rightarrow x_{9,2}$	0.1228
$x_{7,2} \Rightarrow x_{3,2}$	0.0625
...	...
$x_{3,3} \Rightarrow x_{5,2}$	0.0000
$x_{3,3} \Rightarrow x_{5,3}$	0.0000
$x_{4,2} \Rightarrow x_{5,2}$	0.0000
$x_{4,2} \Rightarrow x_{5,3}$	0.0000
$x_{5,2} \Rightarrow x_{3,3}$	0.0000
$x_{5,2} \Rightarrow x_{4,2}$	0.0000
$x_{5,3} \Rightarrow x_{3,3}$	0.0000
$x_{5,3} \Rightarrow x_{4,2}$	0.0000

TABLE III  
CONFIDENCE OF RULES BETWEEN INCIDENT LEVELS.

first blacklisted for 1 minute only and then put back on the white list. In case it is detected misconfigured again, then the blacklist duration is increased to 2 minutes, then to 4 minutes, 16 minutes, etc.

### B. Experiment conditions

Two workflow activities are considered. FIELD-II/pasa consists of 122 invocations of an ultrasonic simulator on an echocardiography 2D data set. It is a data-intensive activity where invocations use from a few seconds to some 15 minutes of CPU time; it transfers 208 MB of input data and outputs about 40 KB of data. Mean-Shift/hs3 has 250 CPU-intensive invocations of an image filtering application. Invocation CPU time ranges from a few minutes up to one hour; input data size is 182 MB and output is less than 1 KB. Files were replicated on two storage sites for both activities.

Two experiments were performed on both workflow activities. Experiment 1 aims at testing that recoverable errors are detected and handled. It is a correct execution where all the input files exist and the application is supposed to run

properly and produce the expected results. Five repetitions were performed for each workflow activity.

Experiment 2 aims at testing that unrecoverable errors are quickly identified and the execution is stopped. Unrecoverable errors were intentionally injected in 3 different runs: in run non-existent inputs, non-existent file paths were used for all the invocations; in application-error, all the file paths existed but input files were corrupted; and in non-existent output, input files were correct but the application did not produce the expected results.

MOTEUR was configured to resubmit failed tasks up to 5 times in all runs of both experiments. For each experiment, a workflow execution using our method (Self-Healing) was compared to a control execution (No-Healing). Executions were launched in production conditions, i.e., without any control of the number of available resources and reliability. Self-Healing and No-Healing were both launched simultaneously to ensure similar grid conditions. Runs were performed along a time period of one week, therefore under different grid conditions. The DIRAC scheduler was configured to equally distribute resources among executions. We used DIRAC v5r12p9 and MOTEUR 0.9.19.

### C. Results and Discussion

Experiment 1: Fig. 6 shows the makespan of FIELD-II/pasa and Mean-Shift/hs3 for the 5 repetitions. The makespan was considerably reduced in all repetitions of both activities. Speed-up values yielded by Self-Healing ranged from 2.6 to 4 for FIELD-II/pasa and from 1.3 to 2.6 for Mean-Shift/hs3.

Table IV shows occurrences of incident levels and associated actions. All recoverable incidents were observed, except  $x_{7,2}$ . For FIELD-II/pasa,  $x_{2,2}$  was the predominant incident due to the data-intensive nature of the application. No blocked activity was detected due to important task replication triggered by low efficiency. For Mean-Shift/hs3, low efficiency and blocked activity almost equally appeared. The total number of replicated tasks for all repetitions was 1,128 for FIELD-II/pasa (i.e. 1.8 task replication per invocation in average) and 644 for Mean-Shift/hs3 (i.e. 0.5 task replication per invocation in average).

Experiment 2: Fig. 7 shows the makespan of FIELD-II/pasa and Mean-Shift/hs3 for the 3 runs

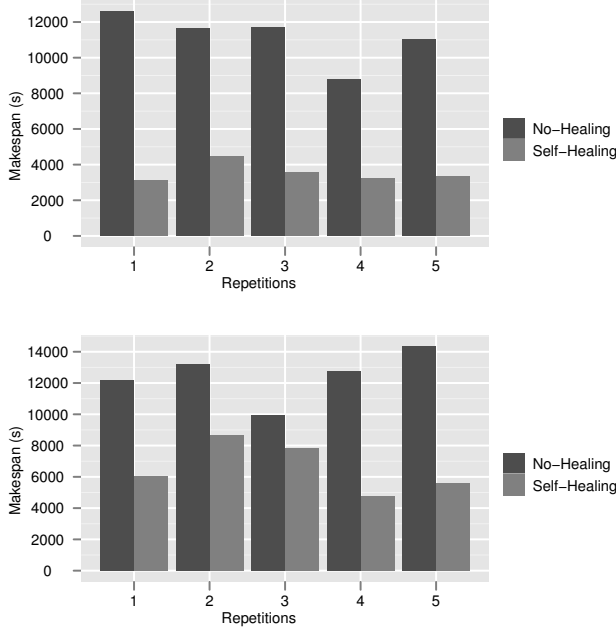


Fig. 6. Experiment 1: execution makespan for FIELD-II/pasa (top) and Mean-Shift/hs3 (bottom).

Activity	Incident level	Occurrence	Actions
FIELD-II/pasa	$x_{2,2}$	262	replicate running tasks replicate input files blacklist site
	$x_{9,2}$	12	
Mean-Shift/hs3	$x_{1,2}$	111	replicate running tasks
	$x_{2,2}$	83	replicate running tasks replicate input files
	$x_{5,2}$	16	replicate files on sites
	$x_{5,3}$	6	blacklist site
	$x_{9,2}$	8	blacklist site

TABLE IV

EXPERIMENT 1: OCCURRENCES OF INCIDENT LEVELS (CUMULATIVE VALUES FOR 5 REPETITIONS).

where unrecoverable errors were introduced. No-Healing was manually stopped after 7 hours to avoid flooding the infrastructure with faulty tasks. In all cases, Self-Healing was able to detect the issue and stop the execution far before No-Healing. It confirms that the healing process is indeed able to identify unrecoverable errors and stop the execution accordingly. As shown on Table V, the number of submitted fault tasks was significantly reduced, which has benefits both to the infrastructure and to the gateway itself.

Run		Number of tasks	
		Self-Healing	No-Healing
application-error	FIELD-II/pasa	196	732
	Mean-Shift/hs3	249	1500
non-existent input	FIELD-II/pasa	293	732
	Mean-Shift/hs3	417	1500
non-existent output	FIELD-II/pasa	287	732
	Mean-Shift/hs3	364	1500

TABLE V

NUMBER OF SUBMITTED FAULTY TASKS.

In average, the experiments used more than 400 nodes from

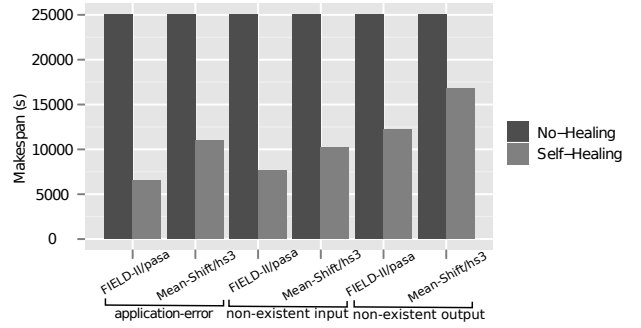


Fig. 7. Experiment 2: makespan of FIELD-II/pasa and Mean-Shift/hs3 for 3 different runs.

36 distinct sites of the production system.

## VII. CONCLUSION

We presented a simple, yet practical method for autonomous detection and handling of operational incidents in workflow activities. No strong assumption is made on the task duration or resource characteristics and incident degrees are measured with metrics that can be computed online. We made the hypothesis that incident degrees were quantified into distinct levels, which we verified using extensive historical information. Incident levels are associated (offline) to action sets ranging from light execution tuning (file/task replication) to radical site blacklisting or activity interruption. Action sets are selected based on the degree of their associated incident level and on confidence of association rules determined from execution history.

This strategy was implemented in the MOTEUR workflow engine and deployed on the European Grid Infrastructure with the DIRAC resource manager. Results show that the proposed method speeds up execution up to a factor of 4 and properly detects unrecoverable errors.

The approach can be extended in several ways. First, other incidents could be added, provided that they can be quantified online by a metric ranging from 0 to 1. Possible candidates are infrastructure service downtimes (e.g. file catalog, storage servers, computing sites) detected by external active monitoring systems such as Nagios [16]. Action sets could also be extended, for instance with actions related to resource provisioning.

Besides, mode detection used for incident quantification could be improved by (i) automated detection (e.g. with Mean-Shift [32]) and (ii) periodical update from execution history. Using the history of actions performed to adjust incident degree could also be envisaged. For instance, incidents for which several actions already have been taken could be considered more critical.

Finally, other gateway components could be targeted with the same approach. Our future work addresses complete workflow executions, taking actions such as pausing workflow executions, detected blocked workflows beyond activities, or allocating resources to users and executions.

## VIII. ACKNOWLEDGMENT

This work is funded by the French National Agency for Research under grant ANR-09-COSI-03 “VIP”. We thank the European Grid Initiative and National Grid Initiatives, in particular France-Grilles, for providing the infrastructure and technical support. We also thank Ting Li and Olivier Bernard for providing optimization use-cases to the Virtual Imaging Platform.

## REFERENCES

- [1] S. Gesing and J. van Hemert, Eds., *Concurrency and Computation: Practice and Experience, Special Issue on International Workshop on Portals for Life-Sciences 2009*, vol. 23, no. 3, march 2011.
- [2] W. Cirne, F. Brasileiro, D. Paranhos, L. Goes, and W. Voorsluys, “On the Efficacy, Efficiency and Emergent Behavior of Task Replication in Large Distributed Systems,” *Parallel Computing*, vol. 33, pp. 213–234, 2007.
- [3] S. Shahand, M. Santcroos, Y. Mohammed, V. Korkhov, A. C. Luyf, A. van Kampen, and S. D. Olabariaga, “Front-ends to Biomedical Data Analysis on Grids,” in *Proceedings of HealthGrid 2011*, Bristol, UK, june 2011.
- [4] P. Kacsuk, “P-GRADE Portal Family for Grid Infrastructures,” *Concurrency and Computation: Practice and Experience*, vol. 23, no. 3, pp. 235–245, 2011.
- [5] R. Ferreira da Silva, S. Camarasu-Pop, B. Grenier, V. Hamar, D. Manset, J. Montagnat, J. Revillard, J. R. Balderrama, A. Tsaregorodtsev, and T. Glatard, “Multi-Infrastructure Workflow Execution for Medical Simulation in the Virtual Imaging Platform,” in *HealthGrid 2011*, Bristol, UK, 2011.
- [6] N. Bard, R. Bolze, E. Caron, F. Desprez, M. Heymann, A. Friedrich, L. Moulinier, N. Nguyen, O. Poch, and T. Toursel, “Déryphon grid - grid resources dedicated to neuromuscular disorders.” in *Stud Health Technol Inform.*, vol. 159, 2010, pp. 124–33.
- [7] D. Malik, J. N. Mordeson, and M. Sen, “On Subsystems of a Fuzzy Finite State Machine,” *Fuzzy Sets and Systems*, vol. 68, no. 1, pp. 83 – 92, 1994.
- [8] J. Kephart and D. Chess, “The vision of autonomic computing,” *Computer*, vol. 36, no. 1, pp. 41 – 50, jan 2003.
- [9] H. Nguyen Van, F. Dang Tran, and J.-M. Menau, “Autonomic virtual resource management for service hosting platforms.” in *Workshop on Software Engineering Challenges in Cloud Computing*, 2009.
- [10] P. Collet, F. Křikava, J. Montagnat, M. Blay-Fornarino, and D. Manset, “Issues and Scenarios for Self-Managing Grid Middleware,” in *Workshop on Grids Meet Autonomic Computing, in association with ICAC’2010*. Washington, DC, USA: ACM, June 2010.
- [11] L. Broto, D. Hagimont, P. Stolf, N. Depalma, and S. Temate, “Autonomic management policy specification in Tune,” in *Proceedings of the 2008 ACM symposium on Applied computing*, New York, NY, USA, 2008, pp. 1658–1663.
- [12] C. Germain-Renaud, A. Cady, P. Gauron, M. Jouvin, C. Loomis, J. Martyniak, J. Nauroy, G. Philippon, and M. Sebag, “The grid observatory,” *IEEE International Symposium on Cluster Computing and the Grid*, pp. 114–123, 2011.
- [13] A. Iosup and D. Epema, “Grid computing workloads,” *Internet Computing, IEEE*, vol. 15, no. 2, pp. 19 –26, march-april 2011.
- [14] D. Lingrand, J. Montagnat, J. Martyniak, and D. Colling, “Analyzing the EGEE production grid workload: application to jobs submission optimization,” in *14th Workshop on Job Scheduling Strategies for Parallel Processing (JSSPP’09)*, Roma, Italy, May 2009, pp. 37–58.
- [15] D. Kondo, B. Javadi, A. Iosup, and D. Epema, “The failure trace archive: Enabling comparative analysis of failures in diverse distributed systems,” in *10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing (CCGrid)*, may 2010, pp. 398 –407.
- [16] E. Imamagic and D. Dobrenic, “Grid Infrastructure Monitoring System Based on Nagios,” in *Proceedings of the 2007 workshop on Grid monitoring*, New York, NY, USA, 2007, pp. 23–28.
- [17] T. Elteto, C. Germain-Renaud, P. Bondon, and M. Sebag, “Towards Non-Stationary Grid Models,” *Journal of Grid Computing*, Dec. 2011.
- [18] X. Zhang, C. Germain-Renaud, and M. Sebag, “Adaptively Detecting Changes in Autonomic Grid Computing,” in *Procs of ACS 2010*, Belgique, Oct. 2010.
- [19] E. Stehle, K. Lynch, M. Shevartalov, C. Torres, and S. Mancoridis, “On the use of computational geometry to detect software faults at runtime,” in *Proceeding of the 7th international conference on Autonomic computing*, New York, NY, USA, 2010, pp. 109–118.
- [20] A. Benoit, L. Marchal, J.-F. Pineau, Y. Robert, and F. Vivien, “Scheduling concurrent bag-of-tasks applications on heterogeneous platforms,” *IEEE Transactions on Computers*, vol. 59, pp. 202–217, 2010.
- [21] C.-C. Hsu, K.-C. Huang, and F.-J. Wang, “Online scheduling of workflow applications in grid environments,” *Future Generation Computer Systems*, vol. 27, no. 6, pp. 860 – 870, 2011.
- [22] H. Casanova, M. Gallet, and F. Vivien, “Non-clairvoyant scheduling of multiple bag-of-tasks applications,” in *Euro-Par 2010 - Parallel Processing*, 2010, vol. 6271, pp. 168–179.
- [23] S. Camarasu-Pop, T. Glatard, J. T. Moscicki, H. Benoit-Cattin, and D. Sarrut, “Dynamic Partitioning of GATE Monte-Carlo Simulations on EGEE,” *Journal of Grid Computing*, vol. 8, no. 2, pp. 241–259, mar 2010.
- [24] H. Casanova, “On the harmfulness of redundant batch requests,” *International Symposium on High-Performance Distributed Computing*, vol. 0, pp. 255–266, 2006.
- [25] W. H. Bell, D. G. Cameron, R. Carvajal-Schiaffino, A. P. Millar, K. Stockinger, and F. Zini, “Evaluation of an economy-based file replication strategy for a data grid,” in *3rd IEEE/ACM International Symposium on Cluster Computing and the Grid*, 2003, p. 661.
- [26] A. H. Elghirani, R. Subrata, and A. Y. Zomaya, “A proactive non-cooperative game-theoretic framework for data replication in data grids,” in *8th IEEE International Symposium on Cluster Computing and the Grid (CCGRID)*, 2008, p. 433.
- [27] J. Rehn, T. Barrass, D. Bonacorsi, J. Hernandez, I. Semeniouk, L. Tuura, and Y. Wu, “Phedex high-throughput data transfer management system.” in *Computing in High Energy Physics, CHEP’2006*, 2006.
- [28] R. Agrawal, T. Imielinski, and A. Swami, “Mining Association Rules between Sets of Items in Large Databases,” 1993, pp. 207–216.
- [29] K. A. De Jong, “An Analysis of the Behavior of a Class of Genetic Adaptive Systems.” Ph.D. dissertation, University of Michigan, Ann Arbor, MI, USA, 1975, aAI7609381.
- [30] T. Glatard, J. Montagnat, D. Lingrand, and X. Pennec, “Flexible and Efficient Workflow Deployment of Data-Intensive Applications on Grids with MOTEUR,” *International Journal of High Performance Computing Applications (IJHPCA)*, vol. 22, no. 3, pp. 347–360, Aug. 2008.
- [31] A. Tsaregorodtsev, N. Brook, A. C. Ramo, P. Charpentier, J. Closier, G. Cowan, R. G. Diaz, E. Lanciotti, Z. Mathe, R. Nandakumar, S. Paterson, V. Romanovsky, R. Santinelli, M. Sapunov, A. C. Smith, M. S. Miguelez, and A. Zhelezov, “DIRAC3. The New Generation of the LHCb Grid Software,” *Journal of Physics: Conference Series*, vol. 219, no. 6, p. 062029, 2009.
- [32] D. Comaniciu and P. Meer, “Mean Shift: A Robust Approach Toward Feature Space Analysis,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 5, pp. 603–619, 2002.